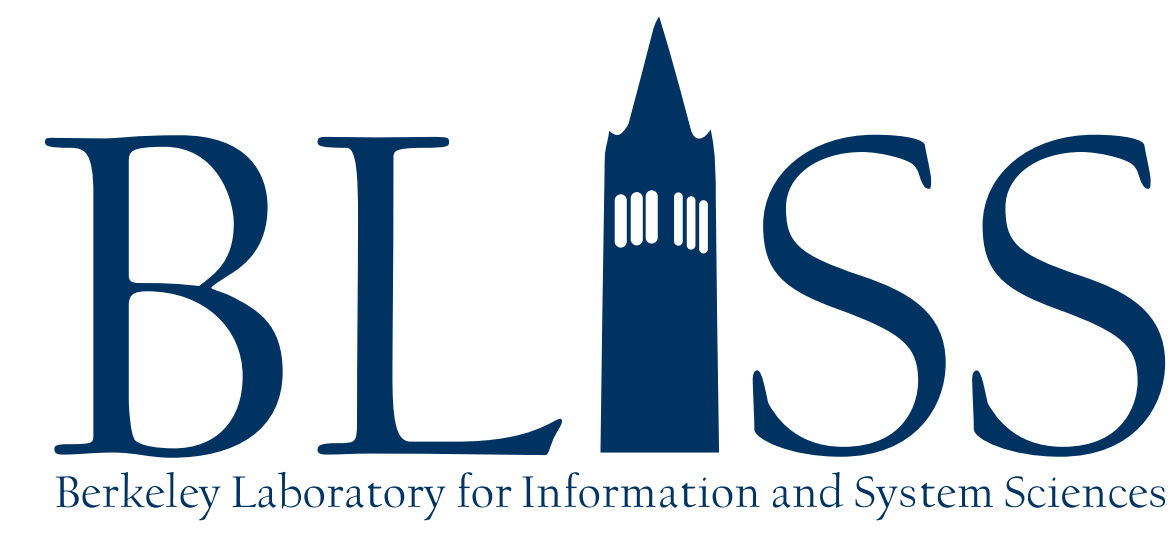




Early stopping for kernel boosting algorithms

Yuting Wei*, Fanny Yang*, Martin Wainwright

Department of Statistics and EECS, University of California, Berkeley



PROBLEM SETTING

- Given **arbitrary regular loss function** $\phi : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$, n fixed covariates x_i and corresponding random $Y_i \sim \mathbb{P}_{x_i}$
- Object of interest: minimizer of population loss function over some **function class \mathcal{F}**

$$\mathcal{L}(f) := \mathbb{E}_{Y_1^n} \left[\frac{1}{n} \sum_{i=1}^n \phi(Y_i, f(x_i)) \right] \text{ and } f^* := \arg \min_{f \in \mathcal{F}} \mathcal{L}(f)$$

- In practice: minimizer of empirical loss function based on observed $\{x_i, Y_i\}_{i=1}^n$

$$\mathcal{L}_n(f) := \frac{1}{n} \sum_{i=1}^n \phi(Y_i, f(x_i)) \text{ and } \hat{f} := \arg \min_{f \in \mathcal{F}} \mathcal{L}_n(f)$$

- If function class large $\mathcal{F} \rightarrow$ risk of **overfitting** to noise!
- Standard way to prevent overfitting: additive penalty function

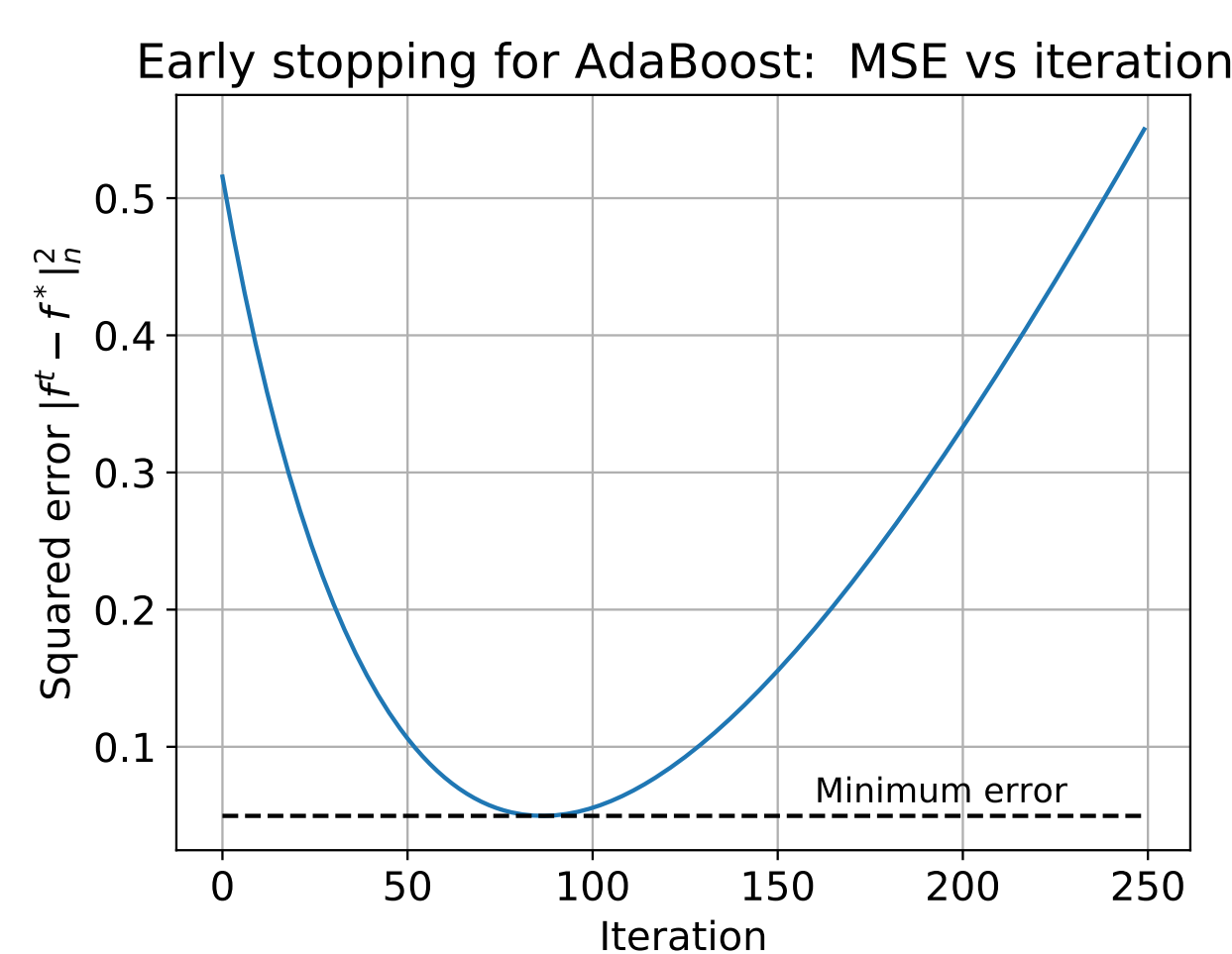
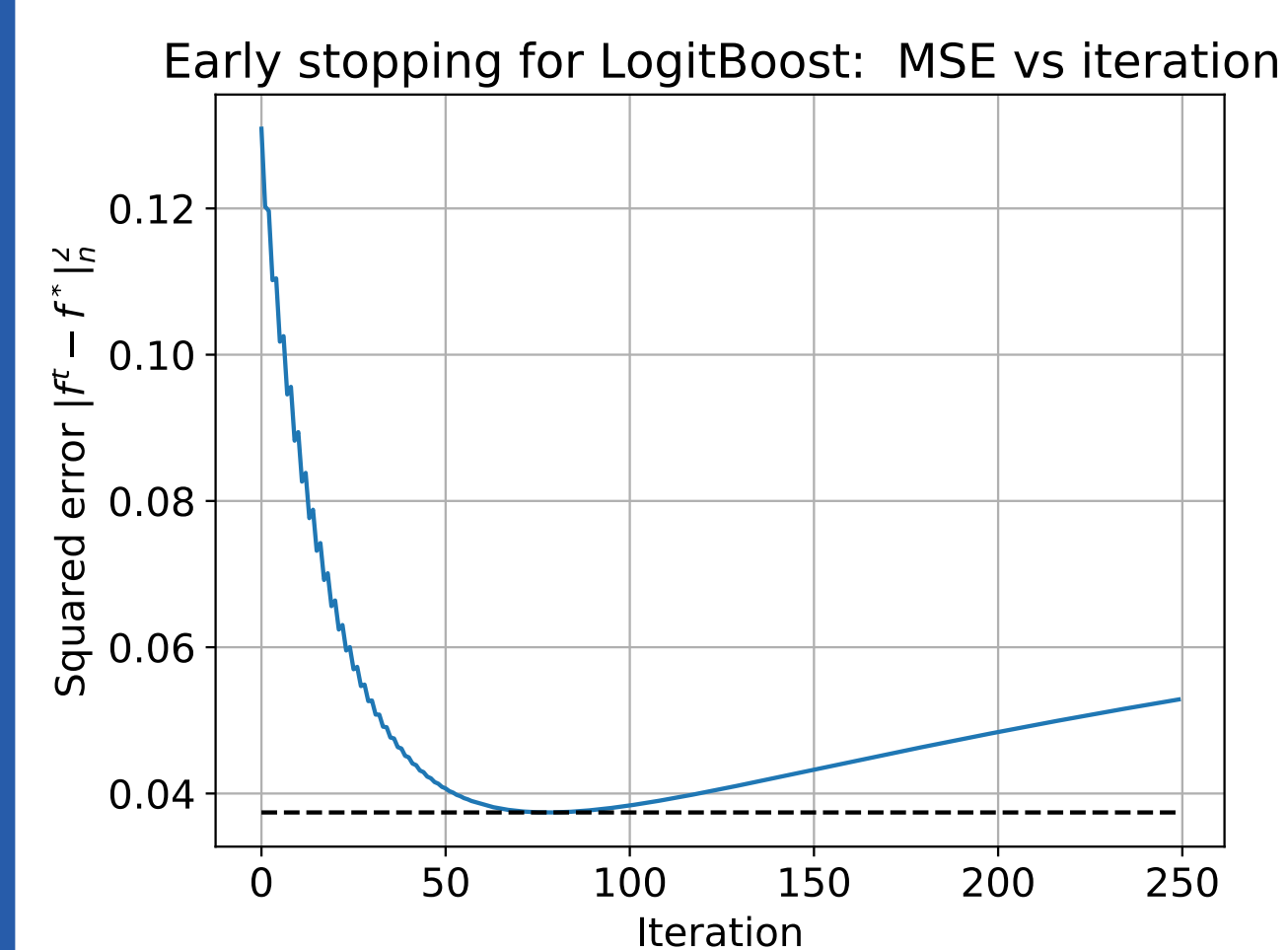
BOOSTING ALGORITHMS

- Based on a sequence of additive updates (weak learners) to improve the fit of a function, see e.g. [1]
- Can be viewed as functional gradient descent steps with updates $f^{t+1} = f^t - \alpha^t g^t$ with $g^t \propto \arg \max_{\|d\|_{\mathcal{F}} \leq 1} \langle \nabla \mathcal{L}_n(f^t), d(x_1^n) \rangle$ (1)

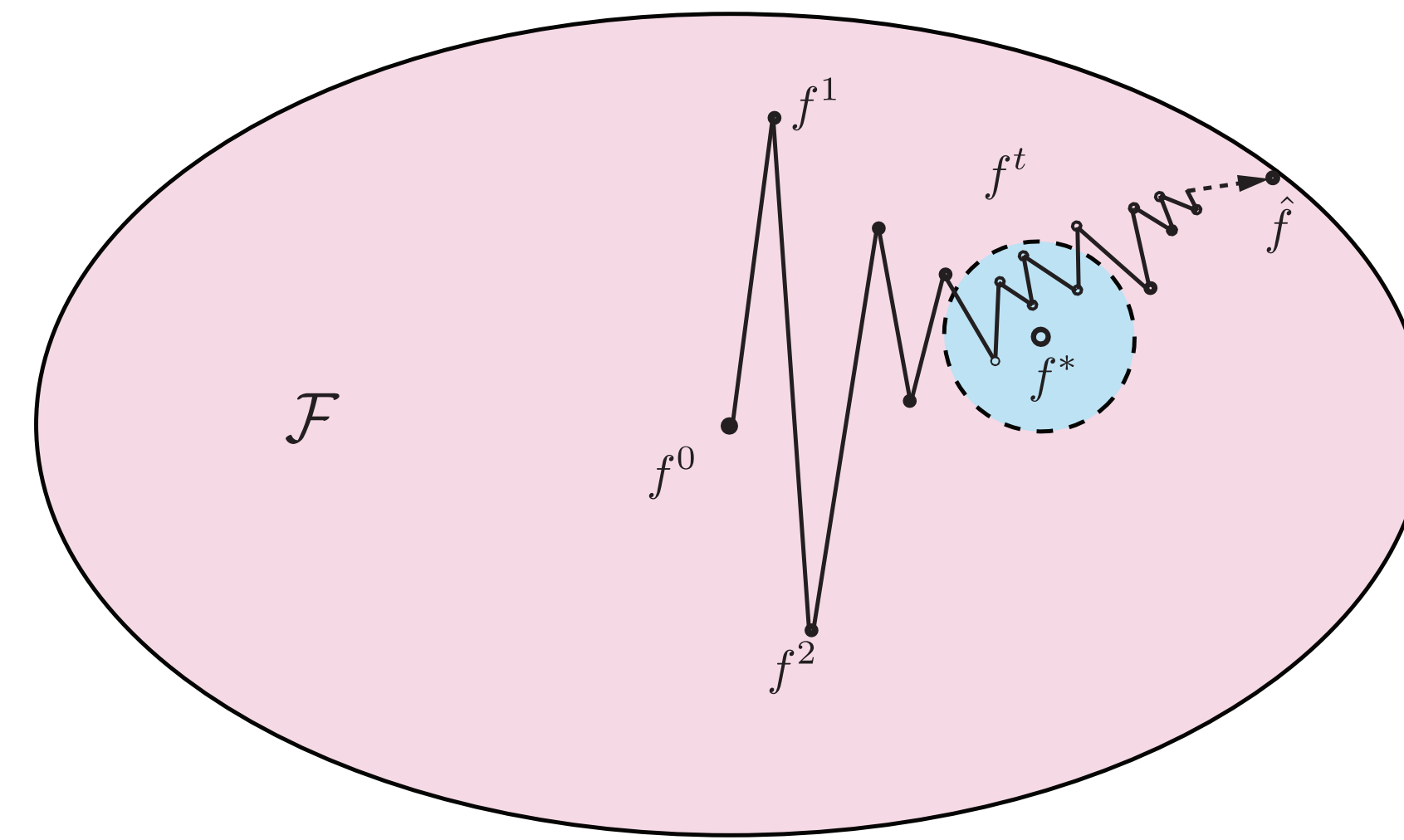
KERNEL BOOSTING

- A positive semidefinite *kernel function* $\mathbb{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ e.g. Gaussian $\mathbb{K}(x, z) = e^{-\frac{(x-z)^2}{2\sigma^2}}$, Sobolev $\mathbb{K}(x, z) = 1 + \min\{x, z\}$ induces a **Reproducing Kernel Hilbert Space** (RKHS) \mathcal{H}
- Optimal $\hat{f} \in \mathcal{H}$ can be represented as $f(\cdot) = \sum_{i=1}^n \omega_i \mathbb{K}(\cdot, x_i)$
- Define kernel matrix K on covariates $\{x_i\}_{i=1}^n$ by $K_{ij} = \mathbb{K}(x_i, x_j)$
- Kernel boosting update** for $\mathcal{F} = \mathcal{H}$ in (1) on vectors

$$f^{t+1}(x_1^n) = f^t(x_1^n) - \alpha n K \nabla \mathcal{L}_n(f^t) \quad (2)$$



OVERFITTING AND EARLY STOPPING



Running until convergence may overfit \rightarrow stop early!

For **least-squares loss**, early stopped boosting (*algorithmic regular.*) and penalized estimators behave similarly statistically, i.e.

$$\|f_{\text{pen}} - f^*\|_n^2 \sim \|f^T - f^*\|_n^2$$

for an appropriate stopping time T , see e.g. [2,3,4].

MAIN CONTRIBUTIONS

- Is there a common principle behind statistical behavior of algorithmic and penalized regularization?

YES! We prove statistical rates for early stopping using the same key quantities as in penalized regularization

- Can we extend to other loss functions?

YES! Our new proof technique allows to extend to a broad class of loss functions (e.g. AdaBoost, LogitBoost ...)

KEY QUANTITIES

- Key quantity I: **Localized Gaussian complexity** $\mathcal{G}_n(\mathcal{E}(\delta, 1)) := \mathbb{E} \left[\sup_{g \in \mathcal{E}(\delta, 1)} \frac{1}{n} \sum_{i=1}^n w_i g(x_i) \right]$, $w_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ where: $\mathcal{E}(\delta, 1) := \left\{ f - g \mid f, g \in \mathcal{H}, \|f - g\|_{\mathcal{H}} \leq 1, \|f - g\|_n \leq \delta \right\}$
- Key quantity II: **Critical radius** δ_n smallest scalar that satisfies

$$\frac{\mathcal{G}_n(\mathcal{E}(\delta, 1))}{\delta} \leq \frac{\delta}{\sigma}$$

REFERENCES

- [1] R. E. Schapire, "The boosting approach to machine learning: An overview," in *Nonlinear estimation and classification*. Springer, 2003, pp. 149–171.
- [2] P. Bühlmann and B. Yu, "Boosting with L^2 loss: Regression and classification," *Journal of American Statistical Association*, vol. 98, pp. 324–340, 2003.
- [3] Y. Yao, L. Rosasco, and A. Caponnetto, "On early stopping in gradient descent learning," *Constructive Approximation*, vol. 26, no. 2, pp. 289–315, 2007.
- [4] G. Raskutti, M. Wainwright, and B. Yu, "Early stopping and non-parametric regression: An optimal data-dependent stopping rule," *Journal of Machine Learning Research*.

MAIN RESULTS

Theorem 1. Given some regular loss function ϕ and the function iterates $\{f^t\}_{t=0}^\infty$ as in (2), for all iterations $T = 0, 1, \dots, \lfloor 1/(8\delta_n^2) \rfloor$, the averaged function estimate \bar{f}^T satisfies with high probability

$$\mathcal{L}(\bar{f}^T) - \mathcal{L}(f^*) \leq C \left(\frac{1}{\alpha T} + \delta_n^2 \right), \quad \text{and} \quad \|\bar{f}^T - f^*\|_n^2 \leq C \left(\frac{1}{\alpha T} + \delta_n^2 \right).$$

Examples for specific kernel spaces:

- γ -exponential decay:** the kernel eigenvalues $\mu_j \leq c_1 \exp(-c_2 j^\gamma)$, when stopped after $T \asymp \frac{n}{\log^{1/\gamma} n}$ steps:

$$\|\bar{f}^T - f^*\|_n^2 \lesssim \frac{\log^{1/\gamma} n}{n}$$

- β -polynomial decay:** the kernel eigenvalues $\mu_j \leq c_1 j^{-2\beta}$, when stopped after $T \asymp n^{2\beta/(2\beta+1)}$ steps:

$$\|\bar{f}^T - f^*\|_n^2 \lesssim n^{-\frac{2\beta}{2\beta+1}}$$

NUMERICAL RESULTS

\mathcal{H} : first order Sobolev space, stop iterates after $T = (2n)^\kappa$

